# Chapter 11 Statistics – (Study) Formula Sheet

## 11.1 – Statistical Studies

Sampling techniques



random sample — equal chance

Systematic sample — Select every 5th

Convenience sample — Easiest (Biased)

Stratified sample — ① split by characteristic ② & Randomly sample within each group (strata)

Cluster sample — ① Mini-populations ② & census each randomly selected group (cluster)
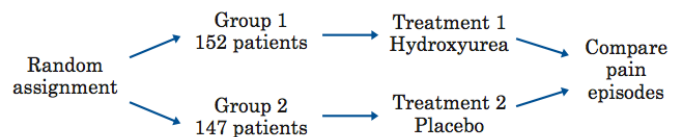
Observational Study vs Experiment

- **Observational study** – <u>Observes</u> existing data.
    - Can reveal association or correlation between variables, but not causation.

- **Experiment** – Generates data to help identify cause-and-effect relationships.
    - <u>Imposes</u> treatments and controls randomly to groups.

Principles of Experimental Design

1. Randomize the control and treatment groups.

2. Control for outside effects on the variable.

3. Replicate the experiment a significant number of times to see meaningful patterns.



Random assignment →
Group 1 152 patients → Treatment 1 Hydroxyurea →
Group 2 147 patients → Treatment 2 Placebo →
Compare pain episodes

## 11.2 – Displaying Data

Frequency Tables

- <u>Summarize datasets</u> by <u>counting</u> the number of observations for each category, distinct value or interval.

| Type of Computer | Frequency | Percent |
|---|---|---|
| Desktop | 11 | 11/50 = 22% |
| Laptop | 23 | 23/50 = 46% |
| Notebook | 9 | 9/50 = 18% |
| Tablet | 7 | 7/50 = 14% |

| Number of Pets | Frequency |
|---|---|
| 1-2 | 7 |
| 3-4 | 3 |
| 5-6 | 3 |
| 7-8 | 2 |

Total = 15

### Examples

a) What percent of observations have between 1 and 4 pets inclusive?

$$\frac{7+3}{15} = \frac{10}{13} = \boxed{66.7\%}$$

Graphical Displays of Data

- Pie charts (categorical data)

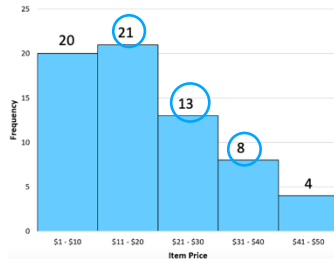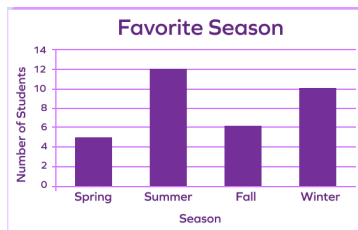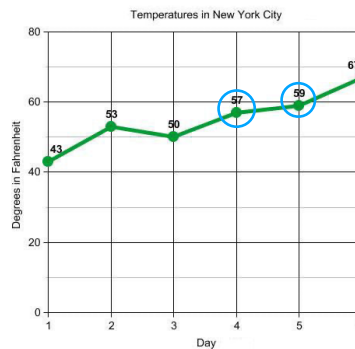    - Compare parts to a whole (slices are proportion of a category).



**Number of Students**
Football 25%
Cricket 17%
Badminton 12%
Hockey 5%
Other 41%

b) What percent of students prefer Football or Hockey?

% football + % Hockey

= 25% + 5%

= $\boxed{30\ \%}$

- Bar graphs (categorical data) and Histograms (numeric data)

    - Height of bar represents amount of data
      in each category (counts or relative frequencies).



- Line graph

    - Shows changes in a numerical
      variable over time.



c) Bar graph – Which season has the highest frequency?

Summer → 12

d) Histogram – How many items cost between $11 and $40 inclusive?

21 + 13 + 8 = 42

e) How many days was the temperature between 55 and 60 °F?

2 days

# 11.3 – Describing and Analyzing Data

Measures of Center

- **Mean** (average) = $\bar{x} = \dfrac{x_1 + x_2 + \cdots + x_n}{n}$

    - NOT resistant → Affected by outliers

- **Median** (middle)

    - The middle value in an ordered list.
    - Resistant → NOT affected by outliers.

- **Mode** (most common)

    - The most frequently occurring value(s).
    - Resistant → NOT affected by outliers.
    - Only measure of center that can be used with categorical data.

Measures of Spread

- **Range** = Max – Min

- **Standard deviation**

    - Measures average distance from the mean.
    - (Don't calculate by hand).

Use calculator to answer these if possible !!!

Example

Dataset: 1, 2, 7, 3, 6, 9, 1, 0, 4, 7

n = 10

a) Find the mean.        ★ **Calc: 1-Var Stats** ☆
                              (Data in L$_1$)

By hand
$\dfrac{1 + 2 + \cdots + 7}{10} = \bar{x} = 4$

b) Find the median.   med = 3.5

0, 1, 1, 2, ③, ④ 6, 7, 7, 9
(3+4)/2 = 3.5

c) Find the mode.

1 + 2 ⟹ occur twice

d) Find the range.

Range = max – min
      = 9 – 0 = 9

e) Find the sample standard deviation.

from calc

S$_x$ = 3.091

↓

Sample

Measures of Relative Position

$P_x$ = some #

Percentile = % LEFT (below)    % Greater Than = 100% - Percentile (RIGHT, above)

- A **percentile** tells you the percent of observations/individuals you are higher than.

- **Quartiles** are specific percentiles.
    - $Q_1$ is the 25th Percentile.
    - $Q_3$ is the 75th Percentile.
    - $Q_2$ is the 50th Percentile = Median.

25%   25% 25%   25%

- **Inner Quartile Range (IQR)**
    - IQR = $Q_3 - Q_1$

1st   2nd 3rd   4th

- **5-number summary**
    - Min, $Q_1$, Med, $Q_3$, Max → Points of a boxplot
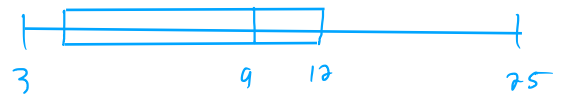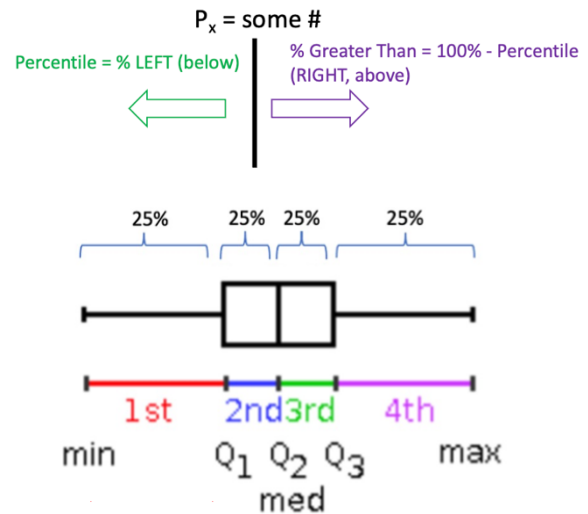
min   $Q_1$ $Q_2$ $Q_3$   max
med

- <u>Example</u>: Calculate the 5-number summary and sketch a boxplot for the following dataset.
    - 12, 3, 4, 7, 21, 3, 9, 8, 10, 11, 25, 11, 13, 4, 5

By hand :   ✗ ✗, ✗⑨ ✗ ✗ ✗,⑨ ✗, ✗ ✗,⑫ ✗,✗, ✗
min =3      $Q_1$ =4       Med = 9       $Q_3$ =12       MAX= 25

3                    9    12                              25

By calc :  1var stats ( $L_1$ = #s)  →   min = 3      $Q_1$ = 12
                                          $Q_{1=}$ 4      MAX = 75
                                          Med = 9

## 11.4 – The Normal Distribution

⭐ Empirical Rule (68 – 95 – 99.7 Rule) ⭐

"step"

<u>68%</u> of the data lies within 1 st dev of the mean.

<u>95%</u> of the data lies within 2 st devs of the mean.

<u>99.7%</u> of the data lies within 3 st devs of the mean.

≈ 99.7%

≈ 95%

≈ 68%

$\mu - 3\sigma$   $\mu - 2\sigma$   $\mu - \sigma$   $\mu$   $\mu + \sigma$   $\mu + 2\sigma$   $\mu + 3\sigma$

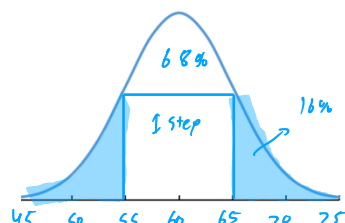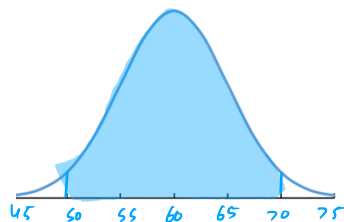- Finding probabilities using the Empirical Rule.

    - Step 1 → **Draw** and **label** curve.
    - Step 2 → **Shade** curve.
    - Step 3 → **Use empirical rule**.

<u>Example</u>
Oak tree heights are normally distributed with mean 60 m and st dev 5 m.

a) Find the percent of trees between 50 m and 70 m tall.

2  steps  ⟹   95 %

45   50   55   60   65   70   75

b) Find the percent of trees greater than 65 m

68%

1 step          16%

Outside = $\frac{\text{Total}}{100\%}$ - $\frac{\text{Inside}}{68\%}$ = 32 %

45   50   55   60   65   70   75

ONLY right = $\frac{32\%}{2}$ = 16%

Finding probabilities based on the normal distribution

- Step 1 → **Standardize** using the **z-score**.

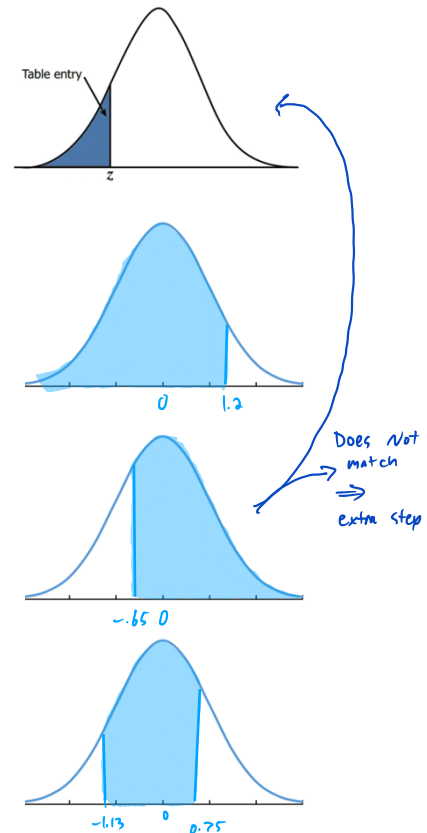  Formula: $z = \dfrac{x-\mu}{\sigma} = \dfrac{obs - mean}{st\ dev}$

  X
  950 970 990 **1010** 1030 1050 1070

  *Standardize*

  Z
  −3 −2 −1 0 +1 +2 +3

  - Ex) $X$ has a normal distribution with mean 10 and st dev 2.
    Find the z-score for $X = 13$.

    $z = \dfrac{13-10}{2} = \boxed{1.5}$

- Step 2 → **Draw**, **label** and **shade** curve.
  - This is how you show your work!!!

  Table entry

- Step 3 → Use '**Standard Normal Distribution**' table to find the probability for Z.
  - Table ALWAYS gives probability LESS THAN Z: P(Z < z).

  - <u>Examples</u>   (How to use table)

  - Left probability = TABLE (Directly)

    $P(Z < 1.2) = \boxed{0.8849}$
    $\quad\quad\quad 1.20$

    0   1.2

    Does Not match
    ⇒
    extra step

  - Right probability = 1 − LEFT

    $P(Z > -0.65) = 1 - P(Z < -0.65) = 1 - 0.2578$
    $\qquad\qquad\qquad = \boxed{0.7422}$

    −.65 0

  - Between probability = LEFT $Z_2$ − LEFT $Z_1$

    $P(-1.13 < Z < 0.75) = P(Z < 0.75) - P(Z < -1.13)$
    $\qquad\quad Z_1 \qquad\quad Z_2$
    $\qquad\qquad\qquad = 0.7734 - 0.1292$
    $\qquad\qquad\qquad = \boxed{0.6442}$

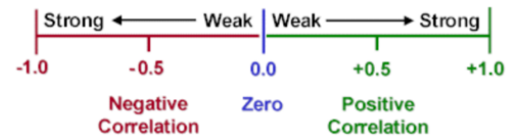    −1.13   0   0.75

## 12.3 – Data Exploration

Scatterplots:

- **Form**: Linear, lurved, or random scatter
- **Direction**: Positive, negative or no association
- **Strength**: Weak, moderate or strong

Perfect Negative Correlation   High Negative Correlation   Moderate Negative Correlation   No Correlation   Moderate Positive Correlation   High Positive Correlation   Perfect Positive Correlation

−1   −0.9   −0.5   0   0.5   0.9   1

Correlation (*r*):

- Interpreting correlation (LINEAR)
    - <u>Sign</u> = Direction
    - <u>Absolute value</u> |*r*| = Strength



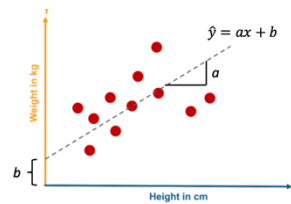- Calculate using calculator
    - **LinReg(ax+b) or 2-Var Stats**      *Show work by writing this*
    - $L_1 = X, L_2 = Y$

Regression:

- Step 1 → Determine if there is a **significant correlation** (**linear relationship**).
    - Compare |*r*| and Critical Value (CV) for *n* (sample size) and significance level $\alpha$.
    - If |*r*| > CV → statistically significant.

| Critical Values of the Pearson Correlation Coefficient | | |
|---|---|---|
| *n* | $\alpha = 0.05$ | $\alpha = 0.01$ |
| 4 | 0.950 | 0.990 |
| 5 | 0.878 | 0.959 |
| 6 | 0.811 | 0.917 |
| 7 | 0.754 | 0.875 |

- Step 2 → Once we have a significant correlation, we can find the **regression line**.
    - $\hat{y} = ax + b$      (get results from correlation calculation)
    - = slope · x + intercept

- Step 3 → Make **predictions** using the regression line.
    - Just plug in the new *X* value to our equation and this will give us the predicted *Y*.



<u>Example</u>

Dataset:

| X | 3 | 5 | 4 | 7 | 6 | 10 |
|---|---|---|---|---|---|---|
| Y | 24 | 40 | 34 | 32 | 17 | 18 |

a) Calculate the correlation *r*.

2-var stats (X,Y)

$r = 0.4205$

OR Linreg(ax+b) → X=$L_1$, Y=$L_2$

b) Determine if *r* is significant for $\alpha = 0.01$.

$n = 6$

$|r| = 0.4205 < 0.917 = CV$

⟹ not significant

c) Suppose we have different regression equation where $\hat{y} = 5x + 2$.

Predict *Y* for *X = 3*:

$\hat{y} = 5(3) + 2 = \boxed{17}$